

---

# Stochastic Collapsed Variational Inference for Hidden Markov Models

---

Pengyu Wang<sup>1</sup> Phil Blunsom<sup>1,2</sup>

<sup>1</sup>Department of Computer Science, University of Oxford

<sup>2</sup>Google DeepMind

{pengyu.wang, phil.blunsom}@cs.ox.ac.uk

## 1 Introduction

Hidden Markov models (HMMs) [1] are popular probabilistic models for modelling sequential data in a variety of fields including natural language processing, speech recognition, weather forecasting, financial prediction and bioinformatics. However, their traditional inference methods such as variational inference (VI) [2] and Markov chain Monte Carlo (MCMC) [3] are not readily scalable to large datasets. For example, one dataset in our experiment consists of 100 million observations.

An important milestone for scaling VI was made by Hoffman et al. [4], who proposed stochastic VI (SVI) that computes cheap gradients based on minibatches of data, updating the model parameters before a complete pass of the full dataset. A recent scalable and more accurate algorithm was proposed by Foulds et al. [5], who applied such stochastic optimization to the collapsed latent Dirichlet allocation (LDA) [6], and their stochastic collapsed variational inference (SCVI) algorithm has been successful in large scale topic modelling.

However, while these recent advances have been studied extensively for topic models that assume a simple bag-of-words data setting [4, 6, 7, 8, 9], there has been little research on whether and how we can apply them in a time dependent data setting. Some research such as SVI for Bayesian time series models [10] and collapsed VI (CVI) for HMMs [11] consider the settings where datasets consist of many independent time series, naturally avoiding to break the sequential dependencies. Perhaps the only true exception is the SVI algorithm for HMMs proposed by Foti et al. [12] in the setting of a single long time series, where the sequential dependencies must be broken.

In this paper, we follow the success of SCVI for LDA [5] and study the SCVI algorithm applied to a single long time series. In a collapsed HMM, we break a long chain into subchains, and we propose a novel sum-product algorithm to update the posteriors of subchains, taking into account their edge transitions due to the sequential dependencies. Our sum product algorithm can be understood as an alternative buffering method to the one in [12]. Our experiments on two discrete datasets show that our SCVI algorithm for HMMs is scalable to very large datasets, memory efficient and significantly more accurate than the existing SVI algorithm.

## 2 Background

A hidden Markov model (HMM) [1] consists of a hidden state sequence  $\mathbf{z} = \{z_t\}_{t=0}^T$  and a corresponding observation sequence  $\mathbf{x} = \{x_t\}_{t=1}^T$ . Let there be  $K$  hidden states. For convenience, we let the start state be 0 and set  $z_0 = 0$ . Let  $\theta$  be the transition matrix where  $\theta_{k,k'} = p(z_t = k' | z_{t-1} = k)$ , and  $\theta_0$  be the initial state distribution where  $\theta_{0,k'} = p(z_1 = k')$ . For  $k = 0, \dots, K$ , we specify the Dirichlet priors with symmetric hyperparameters  $\alpha$  on  $\theta_k$ ,  $\theta_k | \alpha \sim \text{Dir}(\alpha)$  in a Bayesian setting.

A hidden sequence is generated by a Markov process, and each observation is generated conditioned on its hidden state. We have for  $t = 1, \dots, T$ ,

$$z_t | z_{t-1} = k \sim \text{Mult}(\theta_k) \qquad x_t | z_t = k' \sim p(\cdot | \phi_{k'}), \quad (1)$$

where  $\phi_{k'}$  parametrizes the observation likelihood for the hidden state  $k'$ , with  $\phi_{k',w} = p(x_t = w|z_t = k')$ . Without loss of generality, we assume that the observation likelihoods and their conjugate prior take exponential forms. The exponential family is a broad class of probability distributions including multinomial, Gaussian, gamma, Poisson, Dirichlet, Wishart and many others; and there is a conjugate prior distribution for each member in this class. We have for  $k' = 1, \dots, K$ ,

$$p(w|\phi_{k'}) = h_l(w) \exp\{\phi_{k'}^T t(w) - a_l(\phi_{k'})\} \quad (2)$$

$$p(\phi_{k'}|\lambda^\circ) = h_g(\phi_{k'}) \exp\{(\lambda_1^\circ)^T \phi_{k'} + (\lambda_2^\circ)^T (-a_l(\phi_{k'})) - a_g(\lambda^\circ)\}. \quad (3)$$

The base measure  $h$  and log normalizer  $a$  are scalar functions; and the parameter  $\phi_{k'}$  and sufficient statistics  $t$  are vector functions. The subscripts  $l$  and  $g$  represent the local hidden variables and global model parameters, respectively. The dimensionality of the prior hyperparameter  $\lambda^\circ = (\lambda_1^\circ, \lambda_2^\circ)$  is equal to  $\dim(\phi_{k'}) + 1$ .

### 3 Stochastic Collapsed Variational Inference

There is substantial empirical evidence [5, 11, 13] that marginalizing the model parameters is helpful for both accurate and efficient inference. Thus we integrate out the model parameters  $(\theta, \phi)$  and the marginal data likelihood of an HMM is:

$$p(\mathbf{x}, \mathbf{z}) = \prod_{k=0}^K \frac{\Gamma(K\alpha)}{\Gamma(K\alpha + C_{k\cdot})} \prod_{k'=1}^K \frac{\Gamma(\alpha + C_{kk'})}{\Gamma(\alpha)} \prod_{t=1}^T h_l(x_t) \prod_{k'=1}^K \exp\{a_g(\lambda^{k'}) - \{a_g(\lambda^\circ)\}\}. \quad (4)$$

The gamma functions and log normalizers result from the marginalization.  $C_{kk'}$  denotes the transition count from the hidden state  $k$  to  $k'$ ,  $C_{kk'} = \#\{t : z_{t-1} = k, z_t = k'\}$ . dot denotes the summed out column, e.g.,  $C_{k\cdot} = \sum_{k'} C_{kk'}$ .  $\lambda^{k'}$  denotes the posterior hyperparameter for the hidden state  $k'$ ,  $\lambda_1^{k'} = \lambda_1^\circ + \sum_{t=1}^T t(x_t) \delta(z_t = k')$  and  $\lambda_2^{k'} = \lambda_2^\circ + C_{k\cdot}$ , where  $\delta$  is the standard delta function.

Given an observed sequence  $\mathbf{x}$ , the task of Bayesian inference in the collapsed space is to compute the posterior distributions over the hidden sequence,  $p(\mathbf{z}|\mathbf{x})$ . The posteriors over the model parameters can be estimated by taking a variational Bayesian maximization step with our estimated  $q(\mathbf{z})$  [2]. As the exact computation is intractable, we introduce a variational distribution  $q(\mathbf{z})$  in a tractable family and we maximize the evidence lower bound (ELBO) denoted by  $\mathcal{L}(q)$ ,

$$\log p(\mathbf{x}) \geq \mathbb{E}[\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}[\log q(\mathbf{z})] \triangleq \mathcal{L}(q). \quad (5)$$

We consider the tractable family under the generalized mean field assumption [14] in the collapsed space: we break a single long hidden sequence into a set of subchains. We have  $q(\mathbf{z}) = \prod_{n=1}^N q(\mathbf{z}^n)$ . We do not make any further assumptions about the inner structure of each subchain, preserving the inner transition information. It might be worth emphasizing that the time series dependencies in an HMM model are not broken; only the variational posterior is factorized. Therefore, the information can still flow across different subchains via edge transitions.

For notational simplicity, we let each subchain be of the length  $L$  and  $N = \lfloor T/L \rfloor$  be the number of subchains given a long chain. For each hidden subchain  $\mathbf{z}^n = \{z_l^n\}_{l=1}^L$ , we denote the corresponding observed subchain by  $\mathbf{x}^n = \{x_l^n\}_{l=1}^L$ . Combining the work of SCVI for LDA [5] and CVI for HMM [11], we uniformly sample an observation subchain  $\mathbf{x}^n$ , and we derive the posterior update for  $q(\mathbf{z}^n)$  with a zeroth order Taylor approximation [6],

$$q(\mathbf{z}^n) \approx \hat{\theta}_{\cdot, z_1^n} \left( \prod_{l=2}^L \hat{\theta}_{z_{l-1}^n, z_l^n} \right) \hat{\theta}_{z_L^n, \cdot} \left( \prod_{l=1}^L \hat{\phi}_{z_l^n, x_l^n} \right) \quad (6)$$

$$\hat{\theta}_{\cdot, z_1^n} \propto \sum_{z_0^n} q(z_0^n) \left( \mathbb{E}[C_{z_0^n, z_1^n}] + \frac{\alpha}{K q(z_0^n)} \right) \quad (7)$$

$$\hat{\theta}_{z_{l-1}^n, z_l^n} \propto \mathbb{E}[C_{z_{l-1}^n, z_l^n}] + \alpha \quad (8)$$

$$\hat{\theta}_{z_L^n, \cdot} \propto \sum_{z_{L+1}^n} \left( \frac{\mathbb{E}[C_{z_L^n, z_{L+1}^n}] + \frac{\alpha}{K q(z_{L+1}^n)}}{\mathbb{E}[C_{z_L^n, \cdot}] + K\alpha} \right) q(z_{L+1}^n) \quad (9)$$

$$\hat{\phi}_{z_l^n, x_l^n} \propto h(x_l^n) \exp\{a_g(\lambda_1^\circ + t(x_l^n) + \mathbb{E}[t_{z_l^n}(\mathbf{x}, \mathbf{z})], \lambda_2^\circ + 1 + \mathbb{E}[C_{\cdot, z_l^n}])\}, \quad (10)$$

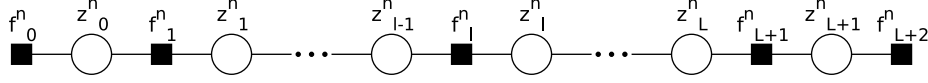


Figure 1: The factor graph of a subchain  $\mathbf{z}^n = \{z_l^n\}_{l=1}^L$  and its guarding variables  $z_0^n$  and  $z_{L+1}^n$ . The emission probabilities have been absorbed into the transition factors.

where  $\mathbb{E}[C_{kk'}] = \sum_{t=1}^T q(z_{t-1} = k, z_t = k')$  denotes the global expected transition count from state  $k$  to  $k'$ , and  $\mathbb{E}[t_{k'}(\mathbf{x}, \mathbf{z})] = \sum_{t=1}^T q(z_t = k') t(x_t)$  denotes the global emission statistics at hidden state  $k'$ . Unlike CVI for HMM [11], we do not need to maintain local statistics and thus our algorithm is memory efficient. We show the algorithmic procedure to infer  $q(\mathbf{z}^n)$  in Section 3.1.

Given  $q(\mathbf{z}^n)$ , we can collect the local transition counts  $\mathbb{E}[C_{kk'}^n]$  and emission statistics  $\mathbb{E}[t_{k'}(\mathbf{x}^n, \mathbf{z}^n)]$  and update the global statistics with an online average weighted by a step size  $\rho_n$ ,

$$\mathbb{E}[C_{kk'}] = (1 - \rho_n) \mathbb{E}[C_{kk'}] + \rho_n T / (L - 1) \mathbb{E}[C_{kk'}^n] \quad (11)$$

$$\mathbb{E}[t_{k'}(\mathbf{x}, \mathbf{z})] = (1 - \rho_n) \mathbb{E}[t_{k'}(\mathbf{x}, \mathbf{z})] + \rho_n N \mathbb{E}[t_{k'}(\mathbf{x}^n, \mathbf{z}^n)]. \quad (12)$$

### 3.1 Modified Forward Backward Algorithm

Given a subchain  $\mathbf{z}^n = \{z_l^n\}_{l=1}^L$ , we denote the hidden variable before it  $z_0^n$  and the hidden variable after it  $z_{L+1}^n$  to be the guarding variables; and we denote  $\hat{\theta}_{z_0^n, z_1^n}$  and  $\hat{\theta}_{z_L^n, z_{L+1}^n}$  to be the edge transitions. In (6), the edge transitions prevent us from applying the standard forward backward algorithm [15] to the HMM parametrized by the surrogate parameters  $\hat{\theta}$  and  $\hat{\phi}$ . Therefore, we propose a modified sum-product algorithm to buffer subchain edges with guarding variables. We start by defining a joint distribution of a subchain and its guarding variables using a factor graph shown in figure 1,

$$q(\mathbf{z}^n, z_0^n, z_{L+1}^n) \propto f_0^n(z_0^n) \left( \prod_{l=1}^{L+1} f_l^n(z_{l-1}^n, z_l^n) \right) f_{L+2}^n(z_{L+1}^n) \quad (13)$$

$$f_0^n(z_0^n) \triangleq q(z_0^n) \quad (14)$$

$$f_1^n(z_0^n, z_1^n) \triangleq \left( \mathbb{E}[C_{z_0^n, z_1^n}] + \frac{\alpha}{Kq(z_0^n)} \right) \hat{\phi}_{z_1^n, x_1^n} \quad (15)$$

$$f_l^n(z_{l-1}^n, z_l^n) \triangleq \hat{\theta}_{z_{l-1}^n, z_l^n} \hat{\phi}_{z_l^n, x_l^n} \quad \text{for } l = 2, \dots, L \quad (16)$$

$$f_{L+1}^n(z_L^n, z_{L+1}^n) \triangleq \frac{\mathbb{E}[C_{z_L^n, z_{L+1}^n}] + \frac{\alpha}{Kq(z_{L+1}^n)}}{\mathbb{E}[C_{z_L^n, \cdot}] + K\alpha} \quad (17)$$

$$f_{L+2}^n(z_{L+1}^n) \triangleq q(z_{L+1}^n). \quad (18)$$

The functions associated with each factor node  $\{f_l^n\}_{l=0}^{L+2}$  are given in (14-18). It is easy to verify that summing over the guarding variables of the joint probability in (13) reduces to  $q(\mathbf{z}^n)$  in (6). Now we can use the sum product algorithm [16] to compute the required marginals of  $q(\mathbf{z}^n)$ . Specifically, we first pick  $f_{L+2}$  as the root node and pass the messages from the leaf node  $f_0$ , and then we pass messages in a reverse direction<sup>1</sup>. We have,

$$u_{f_l^n \rightarrow z_l^n}(z_l^n) = \sum_{z_{l-1}^n} u_{f_{l-1}^n \rightarrow z_{l-1}^n}(z_{l-1}^n) f_l^n(z_{l-1}^n, z_l^n) \quad \text{for } l = 1, \dots, L+1 \quad (19)$$

$$u_{f_{l+1}^n \rightarrow z_l^n}(z_l^n) = \sum_{z_{l+1}^n} u_{f_{l+2}^n \rightarrow z_{l+1}^n}(z_{l+1}^n) f_{l+1}^n(z_l^n, z_{l+1}^n) \quad \text{for } l = L, \dots, 0, \quad (20)$$

where the initial messages are simply the distributions of the two guarding variables,

$$u_{f_0^n \rightarrow z_0^n}(z_0^n) = f_0^n(z_0^n) \quad u_{f_{L+2}^n \rightarrow z_{L+1}^n}(z_{L+1}^n) = f_{L+2}^n(z_{L+1}^n). \quad (21)$$

<sup>1</sup>In both recursions, we have eliminated the messages of the ‘variable node to factor node’ type [17].

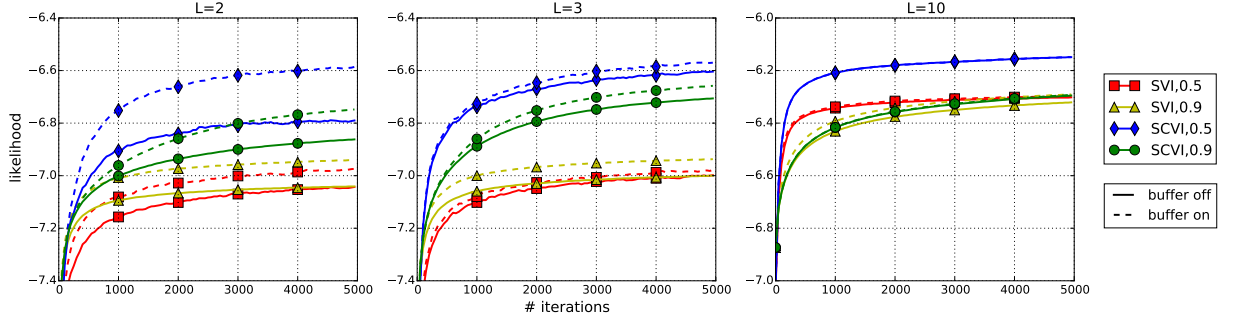


Figure 2: Left and Middle: effect of incorporating buffering methods and performance comparison on WSJ. Right: performance comparison on NYT.

After the messages have been passed in both directions, we compute the required variable marginals  $q(z_l^n)$  and pairwise marginals  $q(z_{l-1}^n, z_l^n)$  by,

$$q(z_l^n) \propto u_{f_l^n \rightarrow z_l^n}(z_l^n) u_{f_{l+1}^n \rightarrow z_l^n}(z_l^n) \quad (22)$$

$$q(z_{l-1}^n, z_l^n) \propto f_l^n(z_{l-1}^n, z_l^n) u_{f_{l-2}^n \rightarrow z_{l-1}^n}(z_{l-1}^n) u_{f_{l+1}^n \rightarrow z_l^n}(z_l^n) \quad (23)$$

The normalization constant can be obtained by normalizing any of these marginals. This completes our algorithm to infer  $q(\mathbf{z}^n)$ .

Our modified sum product algorithm is an alternative buffering method to the one proposed by Foti et al. [12] in their SVI algorithm for a single long time series. A key difference is that we assume the independent subchains and we allow messages to be passed across the borders via local beliefs of the guarding variables in (21), whereas the subchains in the SVI algorithm are naturally correlated. However, the price for preserving the correlation is that they assume the hidden chain is irreducible and aperiodic so that each subchain starts with the initial distribution equal to the stationary distribution of the whole chain. A second superficial difference is that we buffer a subchain by only two guarding variables, whereas Foti et al. buffered a subchain with more observations.

## 4 Experiments

We evaluated the utility of our buffering method and compared the performances of our SCVI algorithm against the SVI algorithm on two synthetic datasets created from the Wall Street Journal (WSJ) and New York Times (NYT). Both corpora are made of sentences, which in turn are sequences of words. For each sentence, the underlying sequence can be understood as a Markov chain of hidden part-of-speech (PoS) tags [18] and words are drawn conditioned on PoS tags, making HMMs natural models. We shuffled both datasets, added special symbols after each sentence to denote the ends and concatenated them. We used the first 1 million words in the concatenated WSJ and 100 million words in the concatenated NYT as our two long time series, respectively. As the evaluation metrics, we used predictive log likelihoods by holding out 5% words of each time series as testing sets.

For both the SVI and our SCVI algorithms: we set the transition and emission priors to be  $\text{Dir}(0.1)$ ; we initialized the global statistics using exponential distributions suggested by Hoffman et al. [4]; we set  $K = 12$  assuming a universal PoS tag set [19]; when buffering was turned off, we set the initial distribution to start a subchain to be the whole chain's stationary distribution. For SVI, when buffering was turned on, we buffered a subchain with 20 words on both sides. We varied the subchain lengths,  $L = 2, 3, 10$  and used minibatches of subchains to reduce the sampling variance. Following Foti et al. [12], we fixed the total length of all subchains in a minibatch  $L \times M = 1000$ , where  $M$  is the minibatch size. Increasing  $L$  means decreasing  $M$  and vice-versa. Also, we varied the forgetting rates  $\kappa = 0.5, 0.9$ , which parametrize the step sizes  $\rho_n = (1 + n)^{-\kappa}$ . Under each of the combined settings, we ran both algorithms for 5000 iterations.

Figure 2 presents the predictive log likelihood results on the WSJ (left and middle) and NYT (right). We see that in most settings our SCVI algorithm outperformed the SVI algorithm by large margins, extending the success of SCVI for LDA [5] to time series data. The only exception is when  $\kappa =$

0.9, both algorithms performed comparably on the NYT. For SCVI, a smaller forgetting rate was preferred, which further promotes the scalability; whereas SVI was less sensitive. When  $L$  is small, there are noticeable improvements using respective buffering methods in both algorithms. For SCVI, we attribute the improvement to the inter subchain communication through guarding variables.

## 5 Conclusion

We have presented a stochastic collapsed variational inference algorithm for HMMs in the setting of a single long time series and an alternative buffering method that modifies the standard forward backward recursions. Our SCVI algorithm is significantly more accurate than the SVI algorithm on two large datasets, and our buffering method is robust against the poor choices of subchain lengths. For future work, we aim to derive the true nature gradients of the ELBO to prove the convergence of our algorithm [20], although we never saw a nonconverging case in our experiments.

## References

- [1] Lawrence R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. pages 267–296, 1990.
- [2] Matthew Beal. *Variational Algorithms for Approximate Bayesian Inference*. PhD thesis, The Gatsby Computational Neuroscience Unit, University College London, 2003.
- [3] L. Scott. Bayesian methods for hidden Markov models: Recursive computing in the 21st century. *Journal of the American Statistical Association*, 97:337–351, 2002.
- [4] Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley. Stochastic variational inference. *J. Mach. Learn. Res.*, 14(1):1303–1347, May 2013.
- [5] James R. Foulds, L. Boyles, C. DuBois, Padhraic Smyth, and Max Welling. Stochastic collapsed variational bayesian inference for latent dirichlet allocation. In *KDD*, 2013.
- [6] Yee Whye Teh, David Newman, and Max Welling. A collapsed variational Bayesian inference algorithm for latent Dirichlet allocation. In *In Advances in Neural Information Processing Systems, volume 19*, 2007.
- [7] Matthew D. Hoffman, David M. Blei, and Francis R. Bach. Online learning for latent dirichlet allocation. In *NIPS*, pages 856–864. Curran Associates, Inc., 2010.
- [8] Chong Wang and David Blei. Truncation-free stochastic variational inference for bayesian nonparametric models. *Advances in Neural Information Processing Systems*, 2012.
- [9] Michael Bryant and Erik B. Sudderth. Truly nonparametric online variational inference for hierarchical dirichlet processes. In F. Pereira, C.J.C. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 2699–2707. Curran Associates, Inc., 2012.
- [10] Matthew Johnson and Alan Willsky. Stochastic variational inference for Bayesian time series models. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1854–1862. JMLR Workshop and Conference Proceedings, 2014.
- [11] Pengyu Wang and Phil Blunsom. Collapsed Variational Bayesian Inference for Hidden Markov Models. In *Proceedings of the 16th International Conference on Artificial Intelligence and Statistics (AISTATS)*, Scottsdale, AZ, USA, 2013.
- [12] Nicholas Foti, Jason Xu, Dillon Laird, and Emily Fox. Stochastic variational inference for hidden Markov models. In *Advances in Neural Information Processing Systems 27*, pages 3599–3607. 2014.
- [13] Arthur Asuncion, Max Welling, Padhraic Smyth, and Yee Whye Teh. On smoothing and inference for topic models. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, Arlington, Virginia, United States, 2009.
- [14] Eric P. Xing, Michael I. Jordan, and Stuart Russell. A generalized mean field algorithm for variational inference in exponential families. In *Proceedings of the Nineteenth Conference on Uncertainty in Artificial Intelligence*, pages 583–591, 2003.
- [15] Leonard E. Baum and Ted Petrie. Statistical inference for probabilistic functions of finite state Markov chains. *Annals of Mathematical Statistics*, 37(6):1554–1563, 1966.
- [16] F. R. Kschischang, B. J. Frey, and H. A. Loeliger. Factor graphs and the sum-product algorithm. *IEEE Trans. Inf. Theor.*, 47(2):498–519, September 2001.
- [17] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.

- [18] Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 1st edition, 2000.
- [19] Slav Petrov, Dipanjan Das, and Ryan T. McDonald. A universal part-of-speech tagset. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC-2012), Istanbul, Turkey, May 23-25, 2012*, pages 2089–2096, 2012.
- [20] Francisco J. R. Ruiz, Neil D. Lawrence, and James Hensman. True natural gradient of collapsed variational bayes. In *NIPS Workshop on Advances in Variational Inference*, Montreal, 2014.